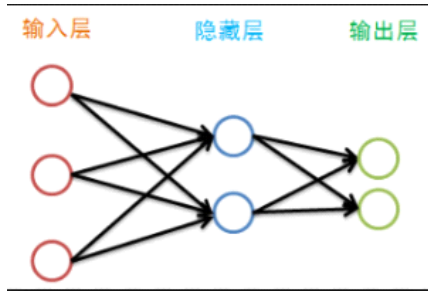


深度学习

2017年12月19日 19:16

简介

1. 结构：中间的隐藏层是无数个机器决定的参数



2. 分类：

standard NN：单层的神经网络

CNN：卷积神经网络，用于图片分析

RNN：序列神经网络，用于音频、语言等序列

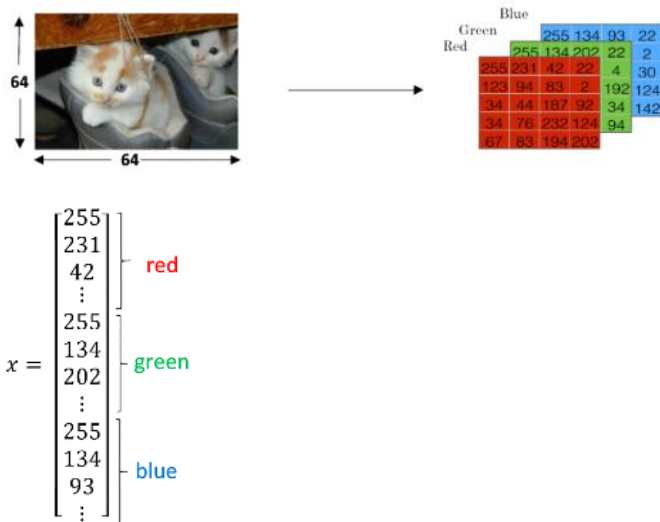
RNNs：多序列的神经网络

3. 发展的原因：

数据量的增长 -> 机器学习 -> 数据量的进一步增长 -> 机器学习不再提升，故采用参数更多的深度学习 -> 对计算速度的要求升高 -> 算法上的改进

数据记录方式

x：图片的基础色可分为3个色：红、绿、蓝，将每个像素点上的亮度均表示出来，得到64*64的3个矩阵，表示为一个64*6*3长度的向量， $x \in R^{n \times x}$



$Y : \{0,1\}$

训练集： $\{(x(1), y(1)), (x(1), y(1)), \dots, (x(m), y(m))\}$, m个样本

矩阵表示：

$$\begin{bmatrix} . & . & . & . \\ . & . & . & . \\ X(1)X(2) & . & X(m) \\ . & . & . & . \\ . & . & . & . \end{bmatrix} \quad [y(1), y(2), \dots, y(m)]$$

$X.shape = (N \times m)$ $Y.shape = (1, m)$

逻辑回归

1. 目的：解决二分类问题

2. 函数：

$$s = \sigma(w^T x + b) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

3. Loss function：

- 正常的损失函数如下，但是该函数是一个非凸的函数，无法利用梯度下降法求解，也没有最优解

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2$$

- 故修改逻辑回归的损失函数如下：

$$L(\hat{y}^{(i)}, y^{(i)}) = -[y^{(i)} \ln \sigma(\hat{y}^{(i)}) + (1 - y^{(i)}) \ln \sigma(1 - \hat{y}^{(i)})]$$

$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

设置成该形式的原因如下:

- 当 $y = 1$ 时, $p(y|x) = \hat{y}$; 当 $y = 0$ 时, $p(y|x) = 1 - \hat{y}$
- 将以上2个式子合为1个式子, 有: $p(y|x) = \hat{y}^y (1 - \hat{y})^{1-y}$
- 两边分别取 \log (原因: \log 是单调函数, 保证 $p(y|x)$ 单调): $\log(p(y|x)) = y \log(\hat{y}) + (1-y) \log(1-\hat{y})$
- 因概率是最大化, 损失函数是最小化, 故前面添加负号: $L = -(y \log(\hat{y}) + (1-y) \log(1-\hat{y}))$

效果:

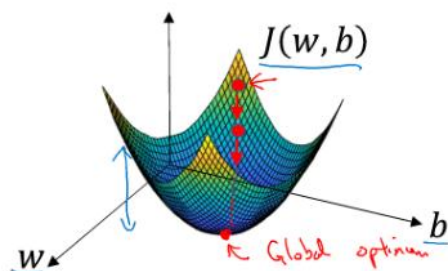
- 当 $y = 1$ 时, $L = -\log(\hat{y})$ 应该要比较小, $\log(\hat{y})$ 应该要比较大, \hat{y} 应该要比较大, \hat{y} 最大就是1
- 当 $y = 0$ 时, $L = -\log(1-\hat{y})$ 应该要比较小, $\log(1-\hat{y})$ 应该要比较大, $1-\hat{y}$ 应该要比较大, \hat{y} 应该要比较小, \hat{y} 最小就是0

4. Cost function: (平均的损失函数)

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})]$$

梯度下降法

1. 希望找到 w 、 b 使得 J 最小

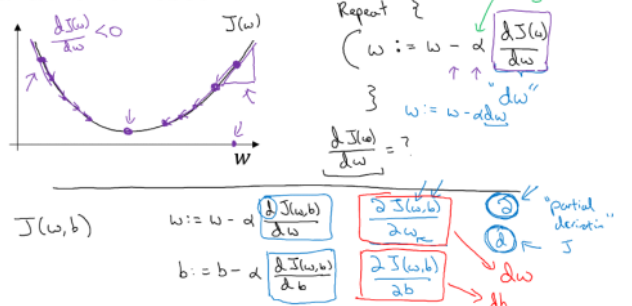


2. w 、 b 给定初始值: 此处因为是凸函数, 初始值给什么样子的最后基本都会收敛到最优值的

3. 迭代:

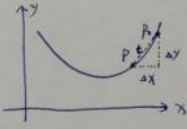
- $w := w - a \frac{\partial J(w, b)}{\partial w}$ $b := b - a \frac{\partial J(w, b)}{\partial b}$
- 当 w 是一个比较大的初始值时, 梯度为正, a 是正数, 新的 w 会比旧的 w 小, 逐步减小直到最优值;
- 当 w 是一个比较小的初始值时, 梯度为负, a 是正数, 新的 w 会比旧的 w 大, 逐步增大直到最优值;

Gradient Descent



4. 为什么导数的方向是梯度下降最快的方向

证明：导数方向是梯度下降最快的方向。



$$\lim_{t \rightarrow 0} \frac{u(P) - u(P_0)}{t}$$

$$= \lim_{t \rightarrow 0} \frac{u(x+\Delta x, y+\Delta y) - u(x, y)}{t}$$

$$= \lim_{t \rightarrow 0} \frac{u(x+\Delta x, y+\Delta y) - u(x, y+\Delta y) + u(x, y+\Delta y) - u(x, y)}{t}$$

$$= \lim_{t \rightarrow 0} \frac{f'_x \cdot \Delta x + f'_y \cdot \Delta y}{t}$$

$$= \lim_{t \rightarrow 0} \frac{f'_x \cdot t \cdot \cos \theta + f'_y \cdot t \cdot \sin \theta}{t}$$

$$= f'_x \cdot \cos \theta + f'_y \cdot \sin \theta$$

$$= \begin{pmatrix} f'_x \\ f'_y \end{pmatrix} \cdot (\cos \theta, \sin \theta)$$

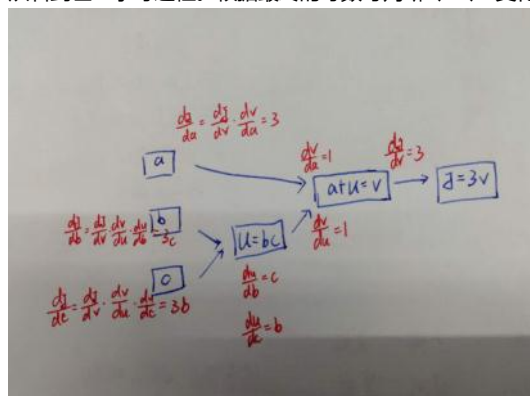
$$= \left| \begin{pmatrix} f'_x \\ f'_y \end{pmatrix} \right| \cdot 1 \cdot \cos \theta$$

当 $\theta = 0$ 时, $\left| \begin{pmatrix} f'_x \\ f'_y \end{pmatrix} \right| \cdot \cos \theta$ 达到最大值。

正向、反向传播介绍

1. 代数例子

- 从左到右：计算过程
- 从右到左：求导过程。根据最终的导数可判断a、b、c变化时对J的影响，从最小化J



2. 逻辑回归例子

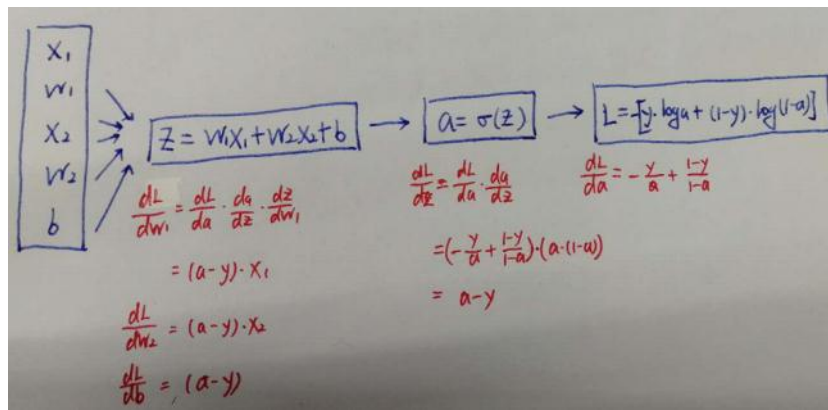
- 理论回顾

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

- 从左到右：计算过程
- 从右到左：求导过程



- 梯度下降参数做自更新

$$w1 := w1 - a \frac{\partial L}{\partial w1} \quad w2 := w2 - a \frac{\partial L}{\partial w2} \quad b := b - a \frac{\partial L}{\partial b}$$

- 扩展到多样本的情况

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))]$$

则分别求出每一个样本的 $\frac{\partial L}{\partial w1}$, $\frac{\partial L}{\partial w2}$, 最后m个样本做平均得到 $\frac{\partial J}{\partial w1}$, $\frac{\partial J}{\partial w2}$

$$\text{梯度下降参数做更新: } w1 := w1 - a \frac{\partial J}{\partial w1} \quad w2 := w2 - a \frac{\partial J}{\partial w2} \quad b := b - a \frac{\partial J}{\partial b}$$

向量化

1. 向量化是一种思想，要时时记得不要用for循环处理数据，多用矩阵
2. 逻辑回归梯度下降的算法实现

参数初始设定

```
J = 0
dw1 = 0
dw2 = 0
db = 0
```

m个样本，n个特征求和

```
for i in range(1,m):
    z[i] = w1*x1[i] + w2*x2[i] + b
    a[i] = sigmod(z[i])
    J += -(y[i]*log(a[i]) + (1-y[i])*log(1-a[i])) # J是m个样本的和
    dz[i] = a[i] - y[i]
    for j in range(1,n):
        dw_j += x[i]_j * dz[i]_j # 计算J对x1的偏导数，是m个样本的
```

求平均值

```
J = J/m
dw1 = dw_1/m
dw2 = dw_2/m
db = db/m
```

参数更新

```
w1 := w1 - adw1
w2 := w2 - adw2
b := b - adb
```

3. 算法优化（利用向量化）

参数初始设定

```
J = 0
dW = np.zeros((n,1))
db = 0
```

m个样本，n个特征求和

```
Z = WT * X + b
A = sigmod(Z)
dZ = A - Y
dW = X * dZ / m
db = np.sum(dZ) / m
```

参数更新

```
W := W - adW
b := b - adb
```

$$\begin{aligned}
 db &= \frac{1}{m} \sum_{i=1}^m dz^{(i)} \\
 &= \frac{1}{m} \text{np.sum}(dz) \\
 dw &= \frac{1}{m} X dz^T \\
 &= \frac{1}{m} \begin{bmatrix} x_1^{(1)} & \dots & x_1^{(m)} \\ \vdots & & \vdots \\ x_n^{(1)} & \dots & x_n^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix} \\
 &= \frac{1}{m} \begin{bmatrix} x_1^{(1)} dz^{(1)} + \dots + x_1^{(m)} dz^{(m)} \\ \vdots \\ x_n^{(1)} dz^{(1)} + \dots + x_n^{(m)} dz^{(m)} \end{bmatrix} \\
 &\quad n \times 1
 \end{aligned}$$

4. 更多向量化的例子

```

u = [1,2,3,4]
v = [6,6,4,1]
np.dot(u,v)
np.exp(v)    # e的v次方
np.log(v)
np.abs(v)
np.maxmun(v,0)
v ** 2

```

5. 广播

- 一个序列 (一个矩阵) + - * / 一个元素
- 一个矩阵 + - * / 一个序列
- `cal = A.sum(axis = 0)` # 求行和
`percentage = A * 100 / cal.reshape(1,4)` # 每个行除以行和

6. Python 技巧

- Python中尽量使用矩阵形式的数据 (5,1) , 不要使用数组形式 (5 ,)

```

In [262]: import numpy as np
          a = np.random.randn(5)
          a.shape

```

Out[262]: (5,)

```

In [263]: a = np.random.randn(5,1)
          a.shape

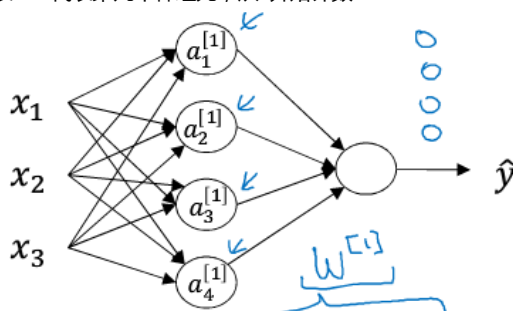
```

Out[263]: (5, 1)

浅层神经网络

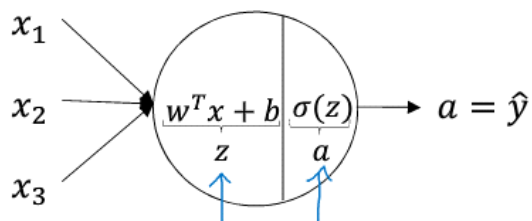
1. 二层神经网络 (输入层为第0层, 隐藏层为第1层, 输出层为第2层)

- 上标 l 代表第几层, 从0开始计数
- 下表 $1 \sim m$ 代表第几个神经元, 从1开始计数



2. 单个样本的二层神经网络计算过程

- 第一层计算过程



$$z = w^T x + b$$

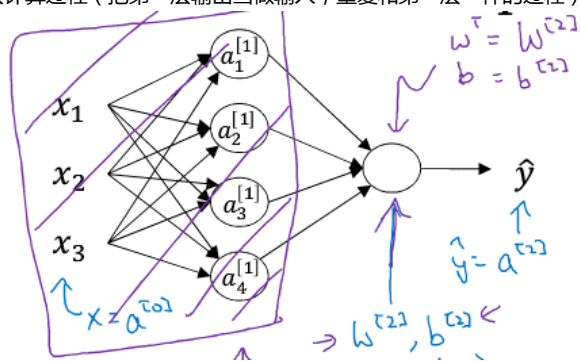
$$a = \sigma(z)$$

$$\begin{aligned} z_1^{[1]} &= w_1^{[1]T} x + b_1^{[1]}, & a_1^{[1]} &= \sigma(z_1^{[1]}) \\ z_2^{[1]} &= w_2^{[1]T} x + b_2^{[1]}, & a_2^{[1]} &= \sigma(z_2^{[1]}) \\ z_3^{[1]} &= w_3^{[1]T} x + b_3^{[1]}, & a_3^{[1]} &= \sigma(z_3^{[1]}) \\ z_4^{[1]} &= w_4^{[1]T} x + b_4^{[1]}, & a_4^{[1]} &= \sigma(z_4^{[1]}) \end{aligned}$$

- 向量化 (由此可完成第一层的计算)

$$\begin{aligned} & \begin{bmatrix} w_{11}^{[1]} \\ w_{12}^{[1]} \\ w_{13}^{[1]} \\ w_{14}^{[1]} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix} = \begin{bmatrix} w_{11}^{[1]}x_1 + b_1^{[1]} \\ w_{12}^{[1]}x_2 + b_2^{[1]} \\ w_{13}^{[1]}x_3 + b_3^{[1]} \\ w_{14}^{[1]}x_4 + b_4^{[1]} \end{bmatrix} = \begin{bmatrix} z_1^{[1]} \\ z_2^{[1]} \\ z_3^{[1]} \\ z_4^{[1]} \end{bmatrix} \\ & \text{尺寸: } (4,1) \quad (4,1) \quad (4,1) \quad (4,1) \quad (4,1) \end{aligned}$$

- 第二层计算过程 (把第一层输出当做输入, 重复和第一层一样的过程)



- 总体向量化

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

(4,1) (4,3) (3,1) (4,1)

$$a^{[1]} = \sigma(z^{[1]})$$

(4,1) (4,1)

$$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$$

(1,1) (1,4) (4,1) (1,1)

$$a^{[2]} = \sigma(z^{[2]})$$

(1,1) (1,1)

3. 多样本的二层神经网络计算过程

单样本过程:

$z_1^{(1)} = W_{01}^{(1)} \cdot x_1^{(1)} + W_{11}^{(1)} \cdot x_2^{(1)} + W_{21}^{(1)} \cdot x_3^{(1)} + b_1^{(1)} \quad a_1^{(1)} = \sigma(z_1^{(1)})$
 $z_2^{(1)} = W_{02}^{(1)} \cdot x_1^{(1)} + W_{12}^{(1)} \cdot x_2^{(1)} + W_{22}^{(1)} \cdot x_3^{(1)} + b_2^{(1)} \quad a_2^{(1)} = \sigma(z_2^{(1)})$
 $z_3^{(1)} = W_{03}^{(1)} \cdot x_1^{(1)} + W_{13}^{(1)} \cdot x_2^{(1)} + W_{23}^{(1)} \cdot x_3^{(1)} + b_3^{(1)} \quad a_3^{(1)} = \sigma(z_3^{(1)})$
 $z_4^{(1)} = W_{04}^{(1)} \cdot x_1^{(1)} + W_{14}^{(1)} \cdot x_2^{(1)} + W_{24}^{(1)} \cdot x_3^{(1)} + b_4^{(1)} \quad a_4^{(1)} = \sigma(z_4^{(1)})$

向量元 1:

$$z_1^{(1)} = W_{01}^{(1)} \cdot x_1^{(1)} + b_1^{(1)} \rightarrow a_1^{(1)} = \sigma(z_1^{(1)})$$

$$z_2^{(1)} = W_{02}^{(1)} \cdot x_1^{(1)} + b_2^{(1)} \rightarrow a_2^{(1)} = \sigma(z_2^{(1)})$$

$$z_3^{(1)} = W_{03}^{(1)} \cdot x_1^{(1)} + b_3^{(1)} \rightarrow a_3^{(1)} = \sigma(z_3^{(1)})$$

$$z_4^{(1)} = W_{04}^{(1)} \cdot x_1^{(1)} + b_4^{(1)} \rightarrow a_4^{(1)} = \sigma(z_4^{(1)})$$

向量元 2:

$$z^{(1)} = \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix} = \begin{bmatrix} -W_{01}^{(1)T} \\ -W_{02}^{(1)T} \\ \dots \\ -W_{04}^{(1)T} \end{bmatrix} \cdot \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{bmatrix} + \begin{bmatrix} b_1^{(1)} \\ b_2^{(1)} \\ b_3^{(1)} \\ b_4^{(1)} \end{bmatrix}$$

$4 \times 1 \quad 4 \times 1 \quad 4 \times 3 \quad 3 \times 1 \quad 3 \times 1$

$$z^{(1)} = \begin{bmatrix} -W_{01}^{(1)T} \\ \vdots \\ -W_{04}^{(1)T} \end{bmatrix} \cdot \begin{bmatrix} z_1^{(1)} \\ z_2^{(1)} \\ z_3^{(1)} \\ z_4^{(1)} \end{bmatrix} + b^{(1)}$$

$1 \times 1 \quad 1 \times 4 \quad 4 \times 1$

总结: 几组样本, 则 W 有几列, 分别要与 X 去乘;
几个输出, 则 W 有几行。

多样本过程:

• 首先, 左侧的神经网络每个样本都有;

• 将 X 按逐列排列, 如下:

$$X = \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & x_1^{(3)} & \dots & x_1^{(n)} \\ x_2^{(1)} & x_2^{(2)} & x_2^{(3)} & \dots & x_2^{(n)} \\ x_3^{(1)} & x_3^{(2)} & x_3^{(3)} & \dots & x_3^{(n)} \end{bmatrix}$$

• 代入左侧已向量化的公式中, 有:

$$z^{(1)} = W^{(1)} \cdot X + b^{(1)}$$

$$A^{(1)} = \sigma(z^{(1)})$$

$$z^{(2)} = W^{(2)} \cdot A^{(1)} + b^{(2)}$$

$$A^{(2)} = \sigma(z^{(2)})$$

其中, $z^{(1)} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & \dots & z^{(1)(n)} \end{bmatrix}$

$$A^{(1)} = \begin{bmatrix} a^{(1)(1)} & a^{(1)(2)} & \dots & a^{(1)(n)} \end{bmatrix}$$

总结: 横向扫描样本, 纵向扫描神经元, 对应不同激活函数。

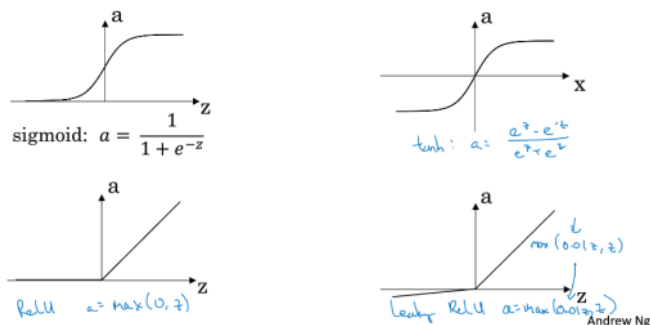
• 推导过程:

$$z^{(1)(1)} = W_{01}^{(1)} x_1^{(1)} + b_1^{(1)} \quad z^{(1)(2)} = W_{02}^{(1)} x_1^{(2)} + b_2^{(1)}$$

$$z^{(1)} = \begin{bmatrix} z^{(1)(1)} & z^{(1)(2)} & \dots & z^{(1)(n)} \end{bmatrix} = W^{(1)} \begin{bmatrix} x_1^{(1)} & x_1^{(2)} & \dots & x_1^{(n)} \\ x_2^{(1)} & x_2^{(2)} & \dots & x_2^{(n)} \\ x_3^{(1)} & x_3^{(2)} & \dots & x_3^{(n)} \end{bmatrix} = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \end{bmatrix}$$

4. 激活函数的选择

- 大背景: 现实生活中, sigmoid函数表现总是不太好, 基本不会使用;
 - 唯一例外: 当二分类的最后一层输出层时可以使用
- 为何要用非线性的激活函数呢?
 - 若激活函数是线性的, 那最终的输出层还只是输入层的线性组合, 不如直接去掉隐藏层, 隐藏层此时不起任何作用。
 - 唯一例外: 回归问题, 在输出层要用恒等激活函数(线性激活函数), 才能输出相应的数值
- 更多的激活函数可以选择

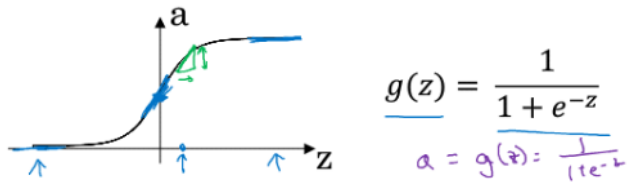


- tanh函数的优点: 将sigmoid函数向下平移, 使有0均值, 相当于正则化;

- relu函数的优点：当x趋于无穷小（或无穷大）时，导数很小，收敛速度很慢，relu不会有这种担忧
- Leaky relu的优点：当x小于0时，为微小的负值

• 激活函数的梯度下降

- sigmoid函数

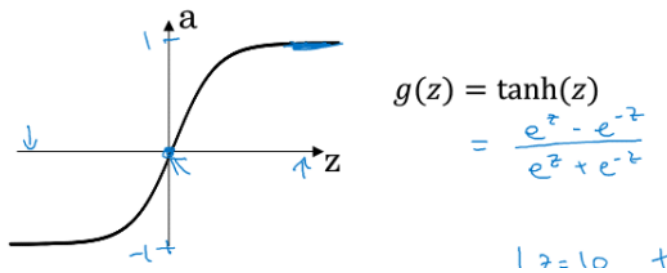


求导：

$$g(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{1 + e^z} = a$$

$$g'(z) = \frac{e^z(1 + e^z) - (e^z)^2}{(1 + e^z)^2} = \frac{e^z}{(1 + e^z)^2} = a(1 - a)$$

- tanh函数



求导：

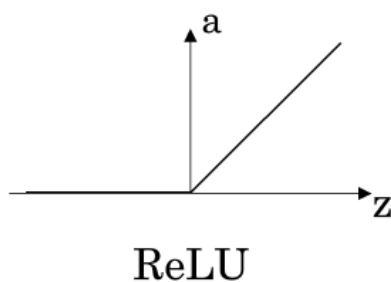
$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} = a$$

$$g'(z) = \left(\frac{e^z + e^{-z} - 2e^{-z}}{e^z + e^{-z}} \right)' = \left(1 - \frac{2e^{-z}}{e^z + e^{-z}} \right)' = \left(1 - \frac{2}{e^{2z} + 1} \right)'$$

$$= -\frac{-e^{2z} \cdot 2}{(e^{2z} + 1)^2} = \frac{4e^{2z}}{(1 + e^{2z})^2}$$

$$= 1 - a^2$$

- relu函数



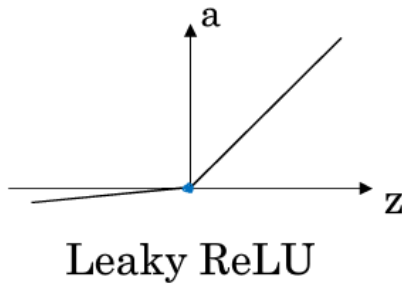
求导:

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1 & z > 0 \\ 0 & z < 0 \\ \text{无定义} & z = 0 \end{cases}$$

但仅一个点处无定义我并care...

- Leaky relu



求导:

$$g(z) = \max(0.01z, z)$$

$$g'(z) = \begin{cases} 0.01 & z < 0 \\ 1 & z > 0 \\ \text{无定义} & z = 0 \end{cases}$$

- 神经网络的梯度下降

- 正向过程

$$\begin{aligned} z^{[1]} &= w^{[1]}x + b^{[1]} \\ A^{[1]} &= g^{[1]}(z^{[1]}) \leftarrow \\ z^{[2]} &= w^{[2]}A^{[1]} + b^{[2]} \\ A^{[2]} &= g^{[2]}(z^{[2]}) = \sigma(z^{[2]}) \end{aligned}$$

- 反向过程

$$\begin{aligned} dz^{[2]} &= a^{[2]} - y \\ dW^{[2]} &= dz^{[2]}a^{[1]T} \\ db^{[2]} &= dz^{[2]} \\ dz^{[1]} &= W^{[2]T}dz^{[2]} * g^{[1]'}(z^{[1]}) \\ dW^{[1]} &= dz^{[1]}x^T \\ db^{[1]} &= dz^{[1]} \end{aligned}$$

$$\begin{aligned} dz^{[2]} &= A^{[2]} - Y \\ dW^{[2]} &= \frac{1}{m} dz^{[2]} A^{[1]T} \\ db^{[2]} &= \frac{1}{m} \text{np.sum}(dz^{[2]}, \text{axis} = 1, \text{keepdims} = \text{True}) \\ dz^{[1]} &= W^{[2]T} dz^{[2]} * g^{[1]'}(z^{[1]}) \\ dW^{[1]} &= \frac{1}{m} dz^{[1]} x^T \\ db^{[1]} &= \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis} = 1, \text{keepdims} = \text{True}) \end{aligned}$$

$J(-) = \frac{1}{m} \sum_{i=1}^m J(\hat{y}_i, y_i)$
 element-wise product

- 证明过程

正向过程:

$$z^{[1]} = w^{[0]} \cdot x + b^{[0]}$$

$$A^{[1]} = g(z^{[1]})$$

$$z^{[2]} = W^{[1]} \cdot A^{[1]} + b^{[1]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

输出层会用 sigmoid 函数。

$$L = y \cdot \log \hat{y} + (1-y) \cdot \log(1-\hat{y})$$

$$= y \cdot \log a + (1-y) \cdot \log(1-a)$$

反向过程:

$$\therefore \frac{\partial L}{\partial a} = \frac{y}{a} - \frac{1-y}{1-a}$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} = \left(\frac{y}{a} - \frac{1-y}{1-a} \right) \cdot \left(\frac{1}{1+e^z} \right)^2 = a - y$$

$$\frac{\partial L}{\partial W} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial W} = (a - y) \cdot x$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial b} = a - y$$

$$\therefore dz^{[2]} = A^{[2]} - y$$

$$dW^{[1]} = \frac{1}{m} (A - y) \cdot A^{[1]T} = \frac{1}{m} dz^{[2]} \cdot A^{[1]T}$$

$$db^{[1]} = \frac{1}{m} \text{np.sum}(A - y, \text{axis}=1, \text{keepdims}=True)$$

$$dA^{[1]} = W^{[1]T} \cdot dz^{[2]} = W^{[1]T} \cdot (A - y)$$

$$dz^{[1]} = dA^{[1]} \cdot g'(z^{[1]})$$

$$dW^{[0]} = \frac{1}{m} dz^{[1]} \cdot x^T$$

$$db^{[0]} = \frac{1}{m} \text{np.sum}(dz^{[1]}, \text{axis}=1, \text{keepdims}=True)$$

$(n^{[2]}, 1)$

$(n^{[2]}, n^{[1]})$

$(n^{[2]}, 1)$

$(n^{[1]}, 1)$

$(n^{[1]}, 1)$

$(n^{[1]}, n^{[0]})$

$(n^{[1]}, 1)$

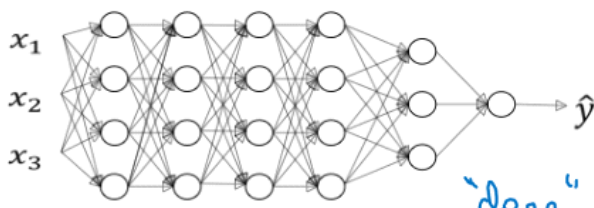
5. 神经网络参数随机化

- 参数一定要随机化, 因为若所有参数均设为0, 那么所有的层和单元都在做一样的运算, 浪费了很多单元
- 随机化的起始参数一定要非常小, 因为大的话, tanh、sigmoid函数的导数接近0, 梯度下降收敛很慢

• 深度神经网络

1. 深度神经网络

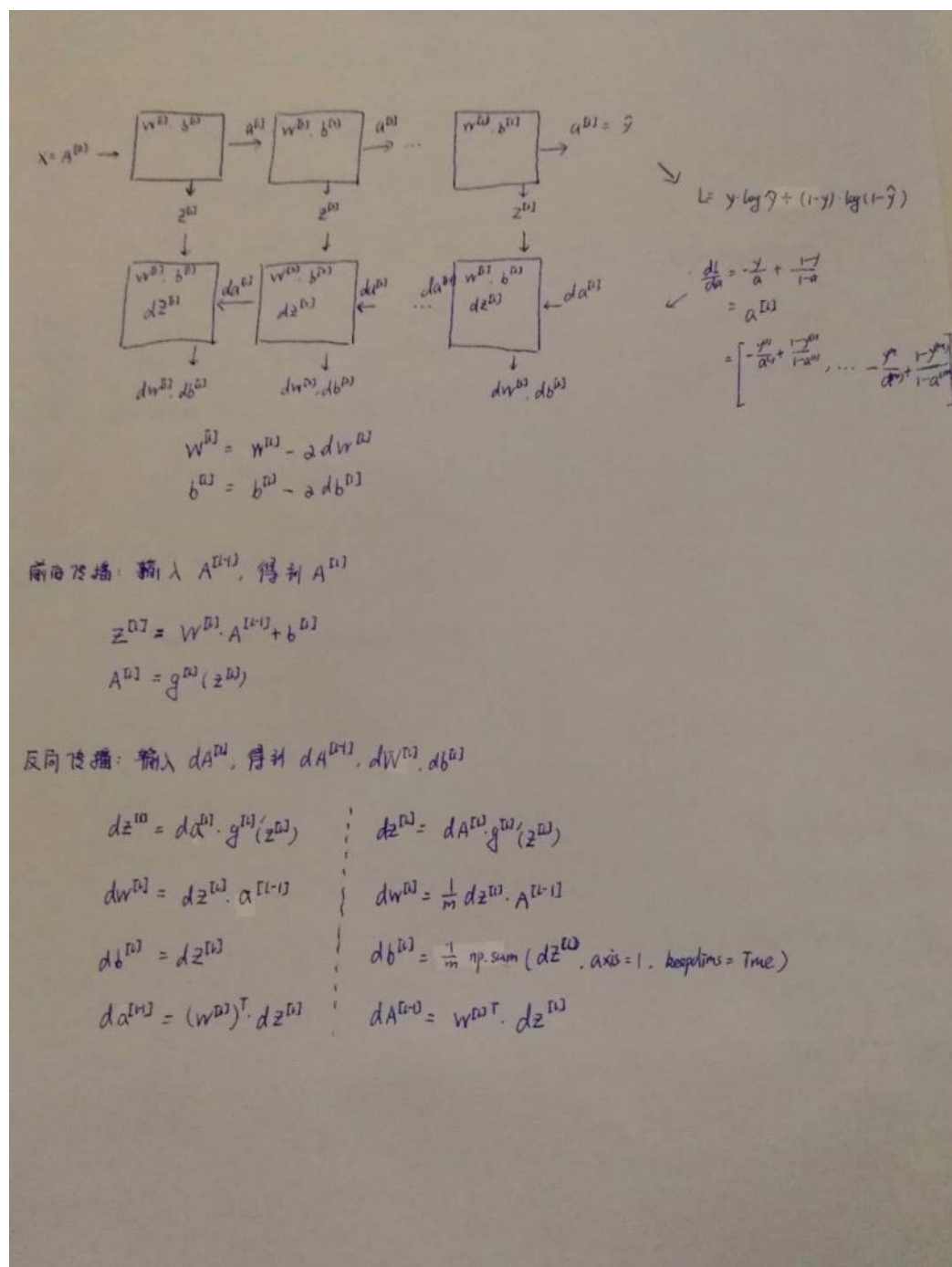
- 层数用L标记



- 为什么要用深度神经网络?

- 神经网络工作时, 浅层的单元通常做的工作是做边缘性的识别, 深层的才做组装, 所以深层的才更有效;
- 深层的神经网络计算量更小

2. 神经网络的工作过程



3. 矩阵维度的介绍

正向过程:

$$z^{[l]} = W^{[l]} \cdot x + b^{[l]}$$

$$A^{[l]} = g^{[l]}(z^{[l]})$$

$$z^{[l+1]} = W^{[l+1]} \cdot A^{[l]} + b^{[l+1]}$$

$$A^{[l+1]} = g^{[l+1]}(z^{[l+1]})$$

$$Y = A^{[L]}$$

第L层的矩阵维度:

$$W^{[L]} : (n^{[L]}, n^{[L-1]}) \quad dw^{[L]} - \text{种}$$

$$b^{[L]} : (n^{[L]}, 1) \quad db^{[L]} - \text{种}$$

$$z^{[L]} : (n^{[L]}, m) \quad dz^{[L]} - \text{种}$$

$$A^{[L]} : (n^{[L]}, m) \quad dA^{[L]} - \text{种}$$

4. 参数和超参数

- 参数: W, b
- 超参数: 会影响和决定参数值的参数
 - 学习率 a
 - 神经网络的层数 L
 - 每层的神经元数 n
 - 每层使用的激活函数 g
 - 算法迭代次数